United States District Court
Northern District of California

1

2

3

4                          UNITED STATES DISTRICT COURT

5                         NORTHERN DISTRICT OF CALIFORNIA

6

7    ABDI NAZEMIAN, et al.,                      Case No.  24-cv-01454-JST   (SK)

8                  Plaintiffs,

9         v.                                     **ORDER REGARDING DISCOVERY**
                                                 **LETTER BRIEF**
10   NVIDIA CORPORATION,
                                                 Regarding Docket No. 162
11                 Defendant.

12        Now before this Court is a dispute about the scope of discovery based on Plaintiffs'

13   requests.  For the reasons set forth below, the Court GRANTS IN PART and DENIES IN PART

14   the request by Plaintiffs.

15        **A.  Background**

16        Plaintiffs Abdi Nazemian, Brian Keene, and Stewart O'Nan (collectively, "Plaintiffs"), on

17   behalf of themselves and all others similarly situated, allege that Defendant NVIDIA Corporation

18   infringed their copyrighted books when it used databases containing that copyrighted material to

19   train its large language models.  (Dkt. No. 1.)  Plaintiffs allege that Defendant announced the

20   availability of four large language models in the "NeMo Megatron series": NeMo Megatron-GPT

21   1.3B, NeMo MegatronGPT 5B, NeMo Megatron-GPT 20B, and NeMo Megatron-T5 3B.  (*Id*.)

22   Defendant publicly disclosed that the NeMo Megatron large language model series was trained on

23   "'The Pile' dataset prepared by EleutherAI."  (*Id*.)  The Pile dataset includes materials from

24   "Books3."  (*Id*.)  Plaintiffs allege that their copyrighted books are included in Books3 and

25   therefore in The Pile.  (*Id*.)  Therefore, Plaintiffs allege that Defendant trained its NeMo Megatron

26   models on Plaintiffs' copyrighted books.  (*Id*.)  Plaintiffs allege direct copyright infringement

27   under 17 U.S.C. § 501.  They purport to represent a class of Plaintiffs from March 8, 2021 to the

28   present, but they intend to amend the class period if they learn that infringement occurred earlier.

1    Plaintiffs seek to represent the following class:

2
3
        All persons or entities domiciled in the United States that own a
        United States copyright in any work that was used as training data for
        the NeMo Megatron large language models during the Class Period.

4    (*Id*.)

5        A short and overly simplified explanation of the process of training a large language model

6    is necessary to understand the dispute.  There are datasets publicly available, and these datasets

7    include large quantities of copyrighted material, used without consent of the owner of the

8    copyrighted material.  Companies such as Defendant then can use those datasets to "train" a large

9    language model to learn patterns of speech, structure, and grammar.  The large language model

10   then learns to predict patterns and then can "understand, generate, and manipulate human

11   language."[1]

12       **B.  Analysis**

13       Plaintiffs seek discovery both about the datasets or libraries that Defendant used to train its

14   large language models and the specific large language models that trained on those datasets.

15   Defendant argues that any discovery should be limited to the large language models that trained on

16   The Pile and that any discovery should also be limited to the four large language models identified

17   in the Complaint.

18       Plaintiffs argue that, because they seek to represent a class of plaintiffs who own any

19   copyrighted material from any source, they are entitled to discovery on all sources and for any

20   large language models in the Nemo Megatron family that trained on those datasets.  In other

21   words, because Plaintiffs do not define the term "family" and do not limit their discovery requests

22   to the four large language models identified in the Complaint, Plaintiffs seek information about

23   any use of any dataset by Defendant.

24       With respect to the datasets, the Court finds that limiting the discovery to the dataset that

25   Plaintiffs know that Defendant used and that contains Plaintiffs' copyrighted books is appropriate.

26   Plaintiffs argue that they are hampered by lack of knowledge because Defendant is the only entity

27
_____

28   [1] For example, this explanation was generated by artificial intelligence in response to a
     question submitted by the Court.

2

that has access to information about its own activities.  This is an inherent problem in any litigation: a plaintiff might not be able to discover wrongdoing by a defendant.  That lack of knowledge does not justify a discovery request without bounds.  Here, Plaintiffs know that Defendant used The Pile, which contains Plaintiffs' copyrighted books, to train its large language models.  Discovery will thus be limited to training on The Pile.

With respect to the limitation on specific large language models, the Court finds that Plaintiffs are entitled to discover whether Defendant trained other large language models on The Pile.  Here, Defendant's proposed limitation to a "family" of large language models is a meaningless distinction because Defendant has not defined "family."  For this reason, the Court will not place an artificial restriction on the large language models that are the subject of Plaintiffs' search.   They may seek information about large language models that use information from The Pile dataset by seeking information about other large language models in the Nemo Megatron family that exist and that trained on The Pile dataset.

Thus, the Court DENIES the motion to compel discovery regarding the use of datasets other than The Pile but GRANTS the motion to compel discovery regarding the use of The Pile to train datasets in the Nemo Megatron family beyond the four large language models specifically named in the Complaint.

**IT IS SO ORDERED**.

Dated: August 21, 2025

_____

SALLIE KIM
United States Magistrate Judge

3